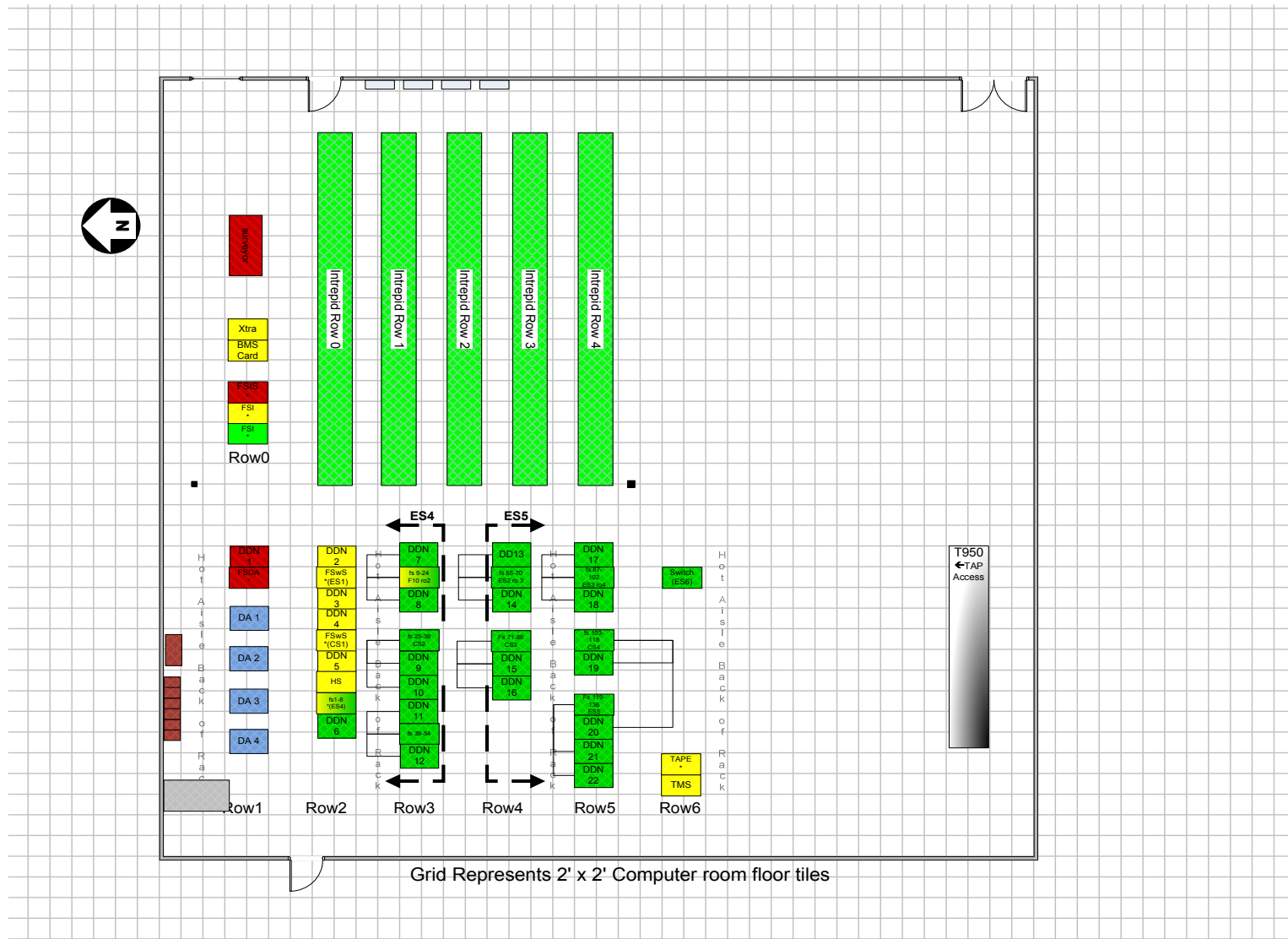


# ALCF Infrastructure Getting Started Workshop

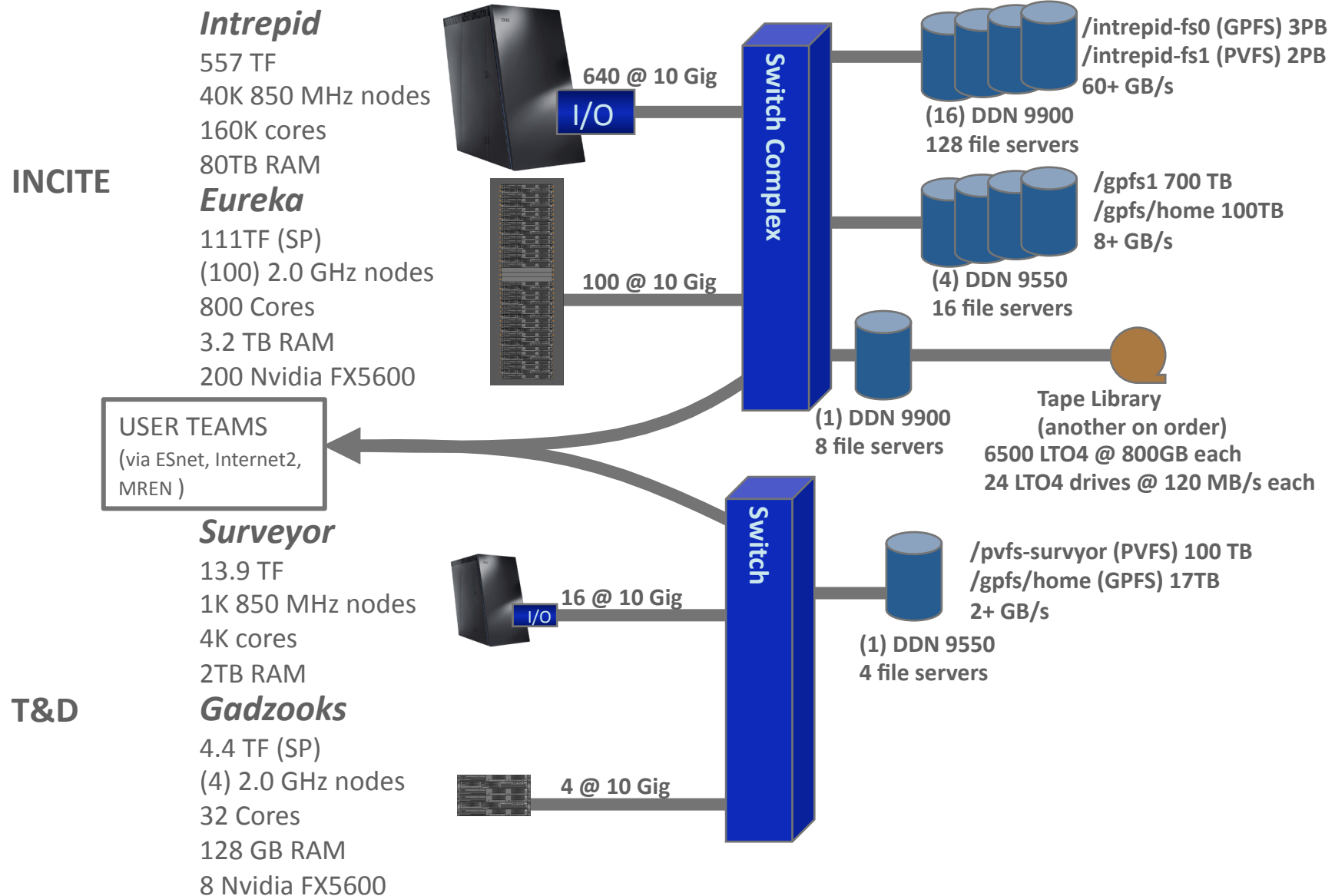
William Scullin

Senior High Performance Computing Systems Administrator  
Leadership Computing Facility

# The World Inside The ISSF



# ALCF Resources



## System Details

Blue Gene /P System	Surveyor	Intrepid (8 racks)
Function	Test & Development	Production INCITE
Login address	surveyor.alcf.anl.gov	intrepid.alcf.anl.gov
Login OS	SLES 10 SP 2	SLES 10 SP 2
Login CPUs	4 PPC970MP @ 2.5 GHz	4 PPC970MP @ 2.5 GHz
Login memory	4GB	4GB
BGP CPU (quad core)	850MHz PPC450fp2	850MHz PPC450fp2
BGP # Nodes / # Cores	1024 / 4096	40,960 / 163,840
BGP Memory	2TB (2GB per node)	80TB (2GB per node)
# I/O nodes (1/64 ratio)	16 @ 10 Gig	640 @ 10 Gig
BGP Compute OS	CNK, ZeptoOS, Plan 9	CNK, ZeptoOS

- Login to ALCF Resources is via ssh with cryptocard authentication
- Round robin DNS will place you on a login node named login[1..n].<machine>.alcf.anl.gov
- Logins are for compilation and job submission only:
  - You may use parallel make, but be gentle
  - Do not do I/O to your home directory



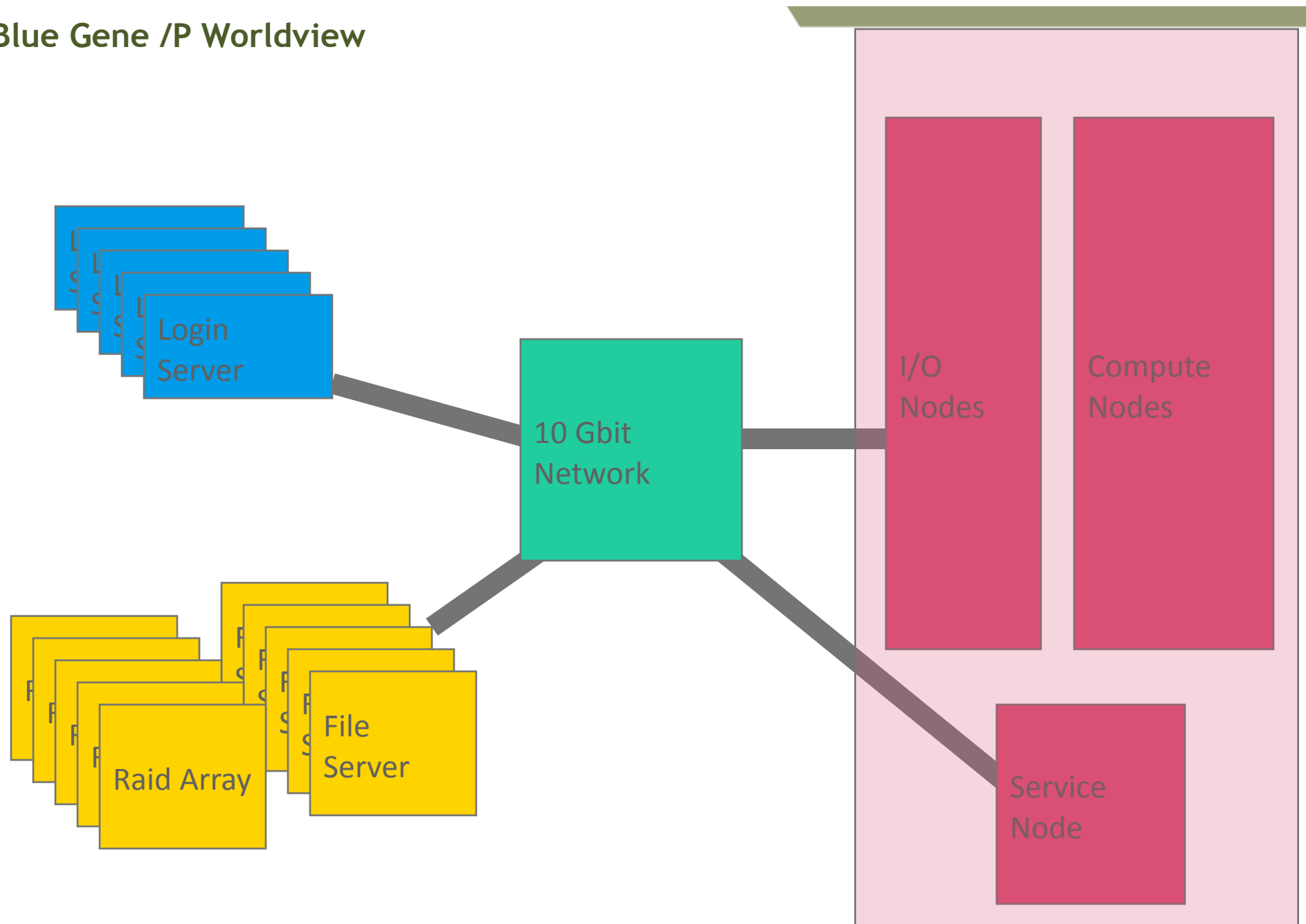
## System Details

Visualization System	Gadzooks	Eureka
Function	Test & Development	Production INCITE
Login access	gadzooks.alcf.anl.gov	eureka.alcf.anl.gov
Node OS	SLES 10 SP 2	SLES 10 SP 2
Node CPU (quad core)	2 Intel Xeon E5405 @ 2.00GHz	2 Intel Xeon E5405 @ 2.00GHz
Nodes / Cores	4 / 32	100 / 800
Memory	128GB (32GB per node)	3.2 TB (32GB per node)
Nvidia FX5600s	8	200
Interconnect	10 GigE (Myrinet)	10 GigE (Myrinet)

- Login to ALCF Resources is via ssh with cryptocard authentication
- Round robin DNS will place you on a login node named login[1..n].<machine>.alcf.anl.gov
- Logins are for compilation and job submission only:
  - You may use parallel make, but be gentle
  - Do not do I/O to your home directory
- Logins are separate from visualization nodes

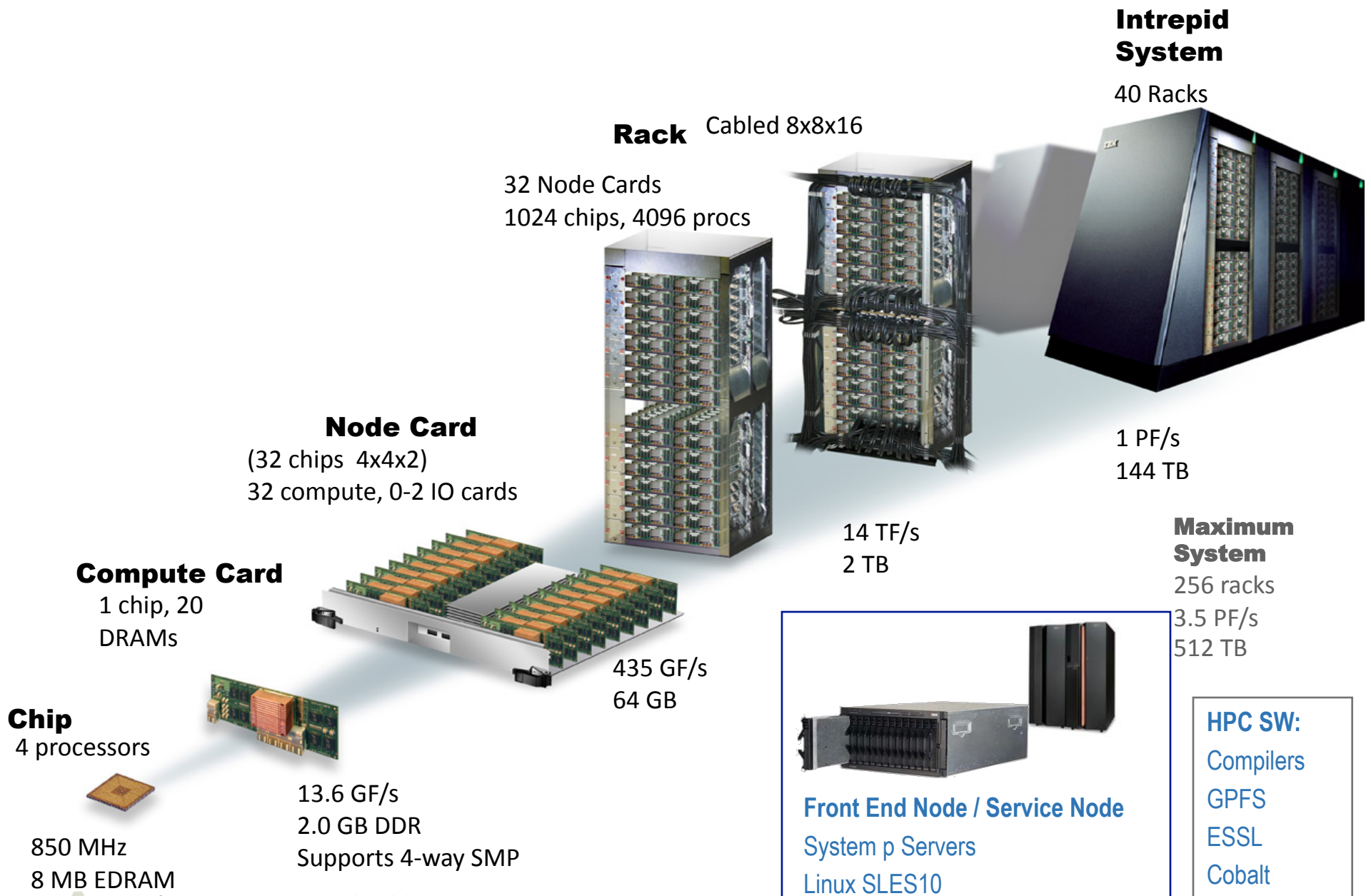


## Blue Gene /P Worldview





# Blue Gene /P Overview

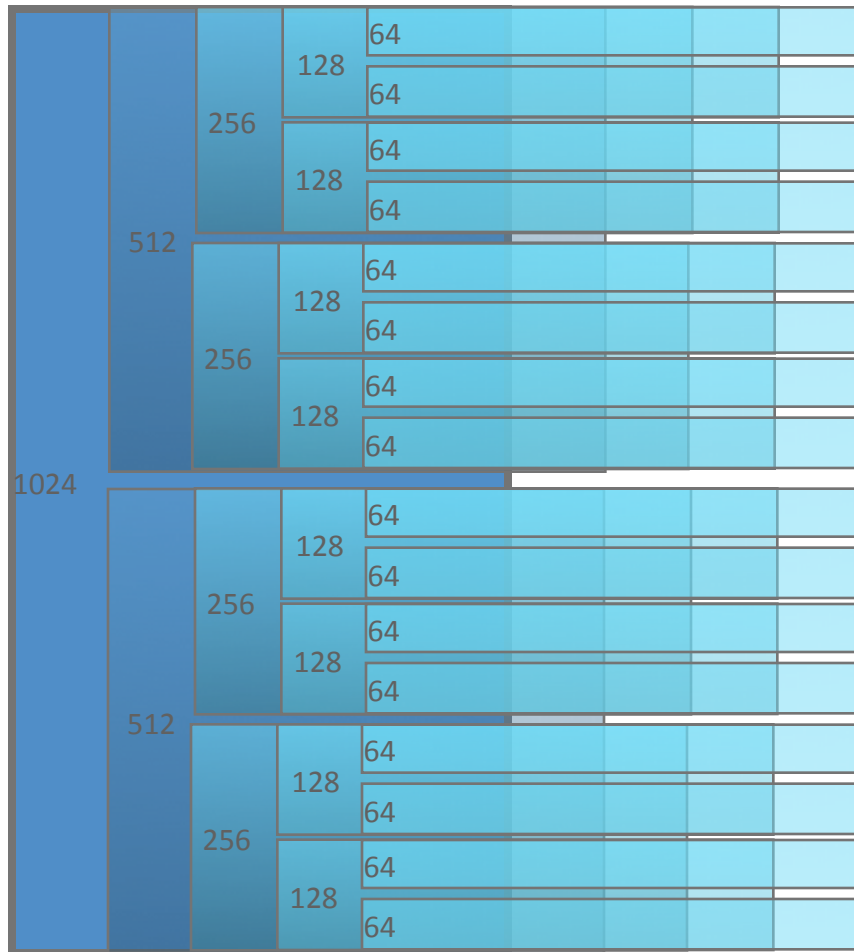


ALCF Infrastructure - Getting Started Workshop



January 27-29, 2010

## Blue Gene Single Rack Partitions (“blocks”)

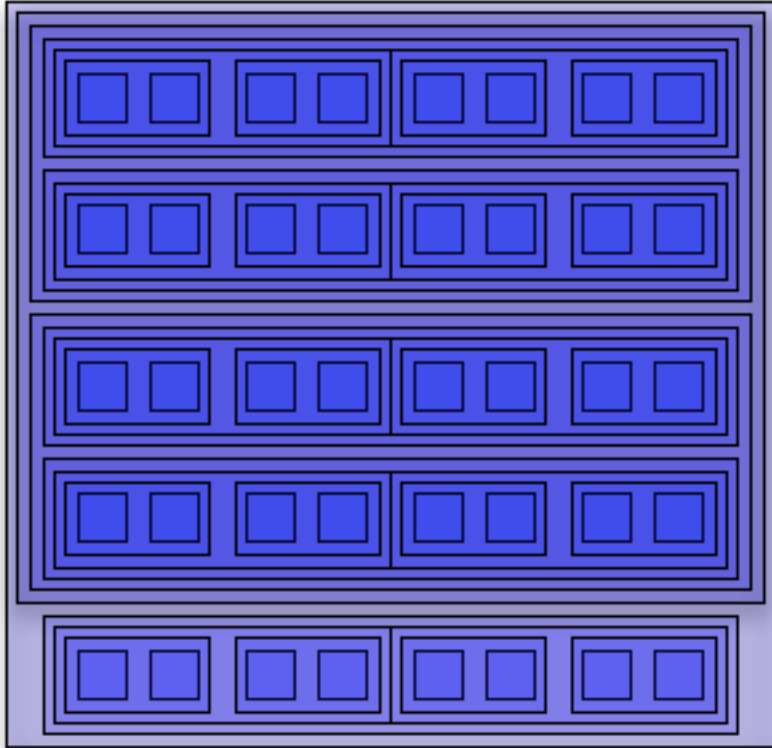


- 1 I/O node for each 64 compute nodes, hardwired to specific set of 64
  - *Minimum partition size of 64 nodes*
- Partition sizes: 64, 128, 256, 512, 1024
  - *Any partition < 512 nodes will get a mesh network layout and not a torus.*
  - *Any partition < 512 nodes will get a non-optimal I/O tree network.*
  - *Do not do performance testing on < 512 nodes*
- Smaller partitions are enclosed inside of larger ones
  - *Not all partitions are available at all times*
  - *Once a job is running on one of the smaller partitions, no jobs can run on the enclosing larger partitions*
- Configuration changes frequently
  - **partlist** shows partition state
- Processes are spread out in a pre-defined mapping, alternate and sophisticated mappings are possible





## Blue Gene Multiple Rack Partitions (“blocks”)



- The following number of large block sizes are possible :
  - 1 40960
  - 1 32768
  - 2 16384
  - 4 8192
  - 9 4096
  - 19 2048
- Not all possible blocks are available at the same time due to wiring dependencies.
- **partlist** will show you if a large free block is busy due to a wiring dependency
- The 40960 node block is generally only available through a reservation
- One rack, R47, is generally reserved for debugging and testing on Intrepid making only the following blocks possible on R4:  
4096, 2 4096
- Mesh partitions are available by reservation only



## Resource Manager and Job Scheduler

- Cobalt - locally developed open source resource manager and scheduler
  - Uses a “cost function” to compute the priority of a job.
- Used on all systems
- Standard commands (qsub, qstat, qdel, qalter)
- Surveyor queues
  - default: Runs the “unicef” cost function; minimize large job starvation while getting good turnaround times
  - Max runtime is 1 hour, no more than 12 jobs running per user, no more than 20 jobs queued per user.
- Intrepid queues
  - prod: Wfp<sup>3</sup> policy; gives priority to larger jobs; will automatically “drain” the machine.
    - Minimum 512 node jobs, max runtime is 12 hours, no more than 5 jobs running per user, no more than 20 jobs queued per user
  - prod-devel: unicef policy (like surveyor)
    - no minimum job size, max runtime is 1 hour, no more than 5 jobs running per user, no more than 20 jobs queued per user

# Reservations

- Reservations
  - Should be the exception not the rule see: <https://wiki.alcf.anl.gov/index.php/Queuing#Reservations> for details
  - Email reservation requests to ***support@alcf.anl.gov***
  - View reservations with **showres**
  - Release reservations with **userres**
- Special reservations
  - R.pm: Preventative maintenance reservation Mondays from 8am to 8am
    - Typically complete in the early evening
  - R.hw\* or R.sw\*: Administrative reservation while addressing hardware or software issues
- This workshop will use:
  - **R.workshop**
  - **R.workshop16**
  - **R.workshop32**

# Allocation Management

- Every user must have at least one Project they are assigned to
- Projects are then given allocations
  - Allocations have an amount, start, and end date and are tracked separately; Charges will cross allocations automatically. The allocation with the earliest end date will be charged first, until it runs out, then the next, and so on
- Charges are based on the partition size, NOT the # of nodes or cores used!
- Reservations are charged for the full time they are active
- Will be managed with clusterbase
  - Use the 'cbank' command (see 'cbank --help')
- Examples:
  - # list all charges against a particular project
    - `cbank -l charge -p <projectname>`
  - # list all active allocations for a particular project
    - `cbank -l allocation -p <projectname>`

## File systems - Intrepid

- Phase II storage: (16) DDN9900 @ 5.5 GB/s each, 128 file servers
  - /intrepid-fs0
    - Intended for very fast parallel IO, program input and output
    - GPFS, 3 PB, 60+ GB/s
    - Not backed up, but you can initiate archive via HPSS
    - Contains
      - /intrepid-fs0/users/\${USER}/scratch
      - /intrepid-fs0/users/\${USER}/persistent
  - /intrepid-fs1
    - Intended for very fast parallel IO, program input and output
    - PVFS, 2 PB, 50+ GB/s
    - Not backed up, but you can initiate archive via HPSS
    - NOTE: Binaries can not be executed from PVFS
    - Contains
      - /intrepid-fs1/users/\${USER}
  - We strongly prefer that users run from /intrepid-fs0 and write to /intrepid-fs0 or /intrepid-fs1. Job I/O to /gpfs/home is viewed as anti-social and is not supported.

## File systems - Intrepid

- Phase I storage: (4) DDN9550 @ 2.2 GB/s each, (8) fs for home, (16) for GPFS
  - /gpfs/home
    - Intended for source code, binaries, etc.. NOT DATA
    - GPFS, 100TB,
    - Backed up via snapshots and tape
  - /gpfs1
    - Intended for fast parallel IO, program input and output
    - GPFS, 700TB, 8+ GB/s
    - Not backed up, but you can initiate archive via HPSS
    - Currently being phased out
- Local storage on login nodes
  - /scratch is available.
    - XFS, 70GB, relatively fast, temporary
    - NOT mounted on BG/P



# File systems - Surveyor

- Phase I storage: (1) DDN9550 @ 2.2 GB/s each, (8) fs for home, (16) for GPFS
  - /gpfs/home
    - Intended for source code, binaries, etc.. NOT DATA
    - GPFS, 15TB,
    - Backed up via snapshots
  - /pvfs-surveyor
    - Intended for fast parallel IO, program input and output
    - PVFS, 88TB
    - Not backed up, CURRENTLY NO TAPE ACCESS
  - We strongly recommend that you avoid writing to /gpfs/home. It is viewed as anti-social and is not supported
- Local storage on login nodes
  - /scratch is available.
    - XFS, 70GB, relatively fast, temporary
    - NOT mounted on BG/P

# Backups and Archival

## ■ Backups

- Snapshots of home directories are done nightly  
~/.snapshot
- home directories are also backed up to tape
  - have not had a single restore request from users
- Data directories will not be backed up

## ■ Archives

- Archive service is available via HPSS
  - HSI is an interactive client
  - HTAR is great for lots of small files
    - Path name is limited to 155 chars in the prefix and 100 bytes for the name (prefix/name)
    - File size is limited to 64 GB.
- GridFTP access to HPSS is available
  - Should be significantly faster

## Getting data in and out

- GridFTP is also available to move data in and out of the site
  - Other site must accept our CA
  - ssh / cryptocard access coming soon.
- Obviously scp is also available.

## Mailing Lists

- For each Blue Gene resource there are two mailing lists.
- Visualization resource related announcements are sent to the mailing list of the associated Blue Gene
- <resource>-users@alcf.anl.gov
  - Mandatory, auto-built from all users with active accounts
  - Important announcements impacting the entire community
    - Security issues
    - Major downtimes
    - Policy changes
    - Long-term news
- <resource>-notify@alcf.anl.gov
  - For the active community
  - Operational status announcements
  - Initially subscribed with account creation
  - Subscribe/unsubscribe as you wish

# Cyber Security

- Argonne computer user agreements
  - Agreed to at account request time
- Standard Argonne computer security rules apply
  - No sharing accounts
    - We WILL know if, for instance, you are letting your grad student use your account.
  - Acceptable use
  - Etc.
- No passwords are allowed for accessing the systems
  - SSH keys used to access Surveyor
  - CRYPTOCARD token required to access Intrepid
    - CRYPTOCARD tokens will work for Surveyor as well
- Data policies are available on the web:
  - <http://www.alcf.anl.gov/support/usingALCF/docs/dataprivacy.php>
  - If you have prohibited data (PII, UNCI, etc.) please contact [support@alcf.anl.gov](mailto:support@alcf.anl.gov)

# Getting Help

## Problems or Questions:

### Check:

- Getting Started: <http://www.alcf.anl.gov/support/usingALCF/usingsystem.php>
- ALCF Wiki: [https://wiki.alcf.anl.gov/index.php/Main\\_Page](https://wiki.alcf.anl.gov/index.php/Main_Page)
- ALCF web pages: <http://www.alcf.anl.gov>
- Intrepid Status: <http://status.alcf.anl.gov/intrepid/activity> (beta)

### Contact:

- e-mail: [support@support.alcf.anl.gov](mailto:support@support.alcf.anl.gov)
- phone: **630-252-3111 (866-508-9181)**
- Your catalyst



Thanks for listening!  
Any questions?